

Datenbank-Spektrum

Das Fachgebiet „Informationssysteme“ am Hasso-Plattner-Institut

--Manuscript Draft--

Manuscript Number:	
Full Title:	Das Fachgebiet „Informationssysteme“ am Hasso-Plattner-Institut
Article Type:	Datenbankgruppen vorgestellt
Corresponding Author:	Felix Naumann Hasso-Plattner-Institut GERMANY
Corresponding Author Secondary Information:	
Corresponding Author's Institution:	Hasso-Plattner-Institut
Corresponding Author's Secondary Institution:	
First Author:	Felix Naumann
First Author Secondary Information:	
Order of Authors:	Felix Naumann Ralf Krestel
Order of Authors Secondary Information:	
Funding Information:	
Abstract:	Das Hasso-Plattner-Institut (HPI) ist ein privat finanziertes Institut an der Universität Potsdam. Stifter ist Professor Hasso Plattner, Mitgründer und Aufsichtsratsvorsitzender des Software-Konzerns SAP. Das Fachgebiet Informationssysteme, das von Prof. Dr. Felix Naumann geleitet wird, beschäftigt sich mit dem effizienten und effektiven Umgang mit heterogenen Daten und Texten. Gegründet wurde das Fachgebiet 2006 und bietet derzeit 12 Doktoranden und circa 15 Masterstudenten eine Forschungsumgebung.
Suggested Reviewers:	

[Click here to view linked References](#)

Das Fachgebiet „Informationssysteme“ am Hasso-Plattner-Institut

Felix Naumann und Ralf Krestel
Hasso-Plattner-Institut
14482 Potsdam
vorname.nachname@hpi.de

Zusammenfassung

Das Hasso-Plattner-Institut (HPI) ist ein privat finanziertes Institut an der Universität Potsdam. Stifter ist Professor Hasso Plattner, Mitgründer und Aufsichtsratsvorsitzender des Software-Konzerns SAP. Das Fachgebiet Informationssysteme, das von Prof. Dr. Felix Naumann geleitet wird, beschäftigt sich mit dem effizienten und effektiven Umgang mit heterogenen Daten und Texten. Gegründet wurde das Fachgebiet 2006 und bietet derzeit 12 Doktoranden und circa 15 Masterstudenten eine Forschungsumgebung.

1. Motivation

Daten sind in großer Fülle vorhanden – man findet sie in vielen verschiedenen Formen von herkömmlichen relationalen oder XML-Datenbanken über semi-strukturierte Daten, oft verlinkt und veröffentlicht als Open Data, bis hin zu Textdaten aus Webdokumenten. Diese Fülle an Daten wird immer größer, und viele Organisationen und Wissenschaftler haben die Vorteile erkannt, diese Daten zu größeren, homogeneren, konsistenteren und saubereren Mengen zu integrieren. Datenintegration vereint unzusammenhängende Quellen in Organisationen; sie bietet den Konsumenten außerdem eine vervollständigte Sicht auf Produktangebote; die Kombination experimenteller Ergebnisse führt zu Gewinnung von neuen wissenschaftlichen Erkenntnissen.

Jedoch ist diese Art von Integration aufgrund der vielfältigen Heterogenitäten schwierig. Syntaktische Heterogenität in Datenformaten, Zugriffsprotokollen und Anfragesprachen ist recht einfach zu lösen, gewöhnlich durch das Bilden von geeigneten quellspezifischen Wrapper-Komponenten. Des Weiteren muss die strukturelle Heterogenität überwunden werden, indem man verschiedene Schemata aneinander ausrichtet. Sogenannte Schema-Matching-Techniken liefern automatisch Ähnlichkeiten und Korrespondenzen entlang der Schemaelemente, während Schema-Mapping-Techniken die Korrespondenzen als eigentliche

Datentransformationen interpretieren. Um schließlich auch die semantische Heterogenität zu überwinden, müssen die unterschiedlichen Bedeutungen von Daten und die ähnliche aber doch unterschiedliche Repräsentation von Echtwelt-Entitäten erkannt werden. Hierfür werden Ähnlichkeitssuche und Datenbereinigungstechniken eingesetzt.

Während wir uns den ersten beiden Herausforderungen in der Vergangenheit gewidmet haben, ist die letzte und wohl auch schwerste, diejenige, auf die wir den Hauptfokus unserer Forschungsarbeit legen. Dieser Fokus lässt sich in drei Hauptforschungsrichtungen aufteilen, welche in den folgenden Kapiteln beschrieben werden: Unser erstes und jüngstes Gebiet ist das Data Profiling, also die Entwicklung von Methoden, um interessante Eigenschaften über unbekannte Datenmengen zu entdecken. Unsere zweite Ausrichtung ist das Gebiet der Datenbereinigung (Data Cleansing), also der Entwicklung von Methoden, um automatisiert Fehler und Unregelmäßigkeiten in Datenbanken zu beheben, insbesondere, um Duplikate zu suchen und zu konsolidieren. Die dritte Forschungsrichtung ist das Text Mining, also das Extrahieren von Informationen aus Textdaten, wie zum Beispiel Wikipedia-Artikel, Tweets und andere Texte aus dem Web.

Soweit es uns möglich ist, versuchen wir unsere Daten und Algorithmen frei zugänglich bereitzustellen:

<http://hpi.de/naumann/projects/repeatability.html>.

2. Data Profiling

„Data Profiling ist eine Folge von Maßnahmen und Prozessen, um Metadaten über eine gegebene Datenmenge zu bestimmen“ [1]. Die Notwendigkeit neue und (noch) unbekannte Daten initial zu untersuchen, wird in vielen Situationen offenkundig, typischerweise in der Vorbereitung auf Folgeaufgaben. Das Profiling umfasst eine breite Palette an Methoden, um effizient Datenmengen zu analysieren. In einem typischen Szenario, welches die

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Fähigkeiten kommerzieller Data-Profiling-Werkzeuge widerspiegelt, werden Tabellen aus relationalen Datenbanken untersucht, um Metadaten zu ermitteln. Dazu gehören Datentyp und typische Wertmuster, Vollständigkeit sowie Einzigartigkeit von Spalten, Schlüssel und Fremdschlüssel, und gelegentlich auch funktionale Abhängigkeiten sowie Assoziationsregeln. Zusätzlich hat die Forschung (sowohl unsere eigene als auch andere) Ansätze für die Entdeckung weiterer Metadaten hervorgebracht, wie z.B. die der Entdeckung von Inklusionsabhängigkeiten oder bedingten funktionalen Abhängigkeiten. Es gibt eine Vielzahl an konkreten Nutzungsbeispielen für die Ergebnisse des Profiling:

- **Anfrageoptimierung:** Anzahlen und Histogramme für die Selektivitätsschätzung, Abhängigkeiten für Anfragevereinfachung
- **Datenbereinigung:** Muster- und Abhängigkeitserkennung, um dann Verstöße festzustellen
- **Datenintegration:** Inklusionsabhängigkeiten über Datenbanken hinweg, um Daten anzureichern und Join-Pfade zu finden
- **Datenanalyse:** Datenaufbereitung und erste Einblicke gewinnen
- **Rekonstruktion von Datenbanken:** Entdeckung von Fremdschlüsseln, um Schemata zu verstehen und deren wichtigsten Komponenten zu identifizieren

Unser Überblicksartikel zeigt den Fortschritt der Community in diesem Bereich [1]. Data Profiling erfreut sich einer immer größeren Aufmerksamkeit, da Forscher und Praktiker zunehmend erkennen, dass das bloße Zusammentragen von Daten in einen großen See (Data Lake) nicht ausreichend ist. „Wenn wir bloß einen Haufen an Datensätzen in einer Ablage sammeln, ist es unwahrscheinlich, dass jemals jemand in der Lage sein wird, diese aufzufinden, geschweige denn diese nochmalig zu verwenden. Mit angemessenen Metadaten besteht jedoch Hoffnung [...]“ [4].

2.1 Profiling relationaler Daten

Abgesehen von den weniger komplexen Aufgaben, wie das Ermitteln der Anzahl unterschiedlicher Werte in einer Spalte, betrachtet Data Profiling komplexe Aufgaben, wie zum Beispiel alle Abhängigkeiten in einer großen Datenmenge zu ermitteln. Wir und auch andere Forschungsgruppen haben verschiedenartige Methoden zur effizienten

Entdeckung aller minimalen funktionalen Abhängigkeiten, Inklusionsabhängigkeiten, eindeutiger Spaltenkombinationen und Ordnungsabhängigkeiten entwickelt. Weitere Abhängigkeiten sollen folgen, beispielsweise Join-Abhängigkeiten, Matching-Abhängigkeiten oder sogenannte Denial Constraints. Anstatt nun jede Profiling-Aufgabe detailliert zu beschreiben, heben wir die allgemeinen Schwierigkeiten hervor, die Data Profiling besonders herausfordernd, aber auch spannend machen:

Schemagröße: Da Abhängigkeiten in jeglichen Spalten und auch Spaltenkombinationen auftreten können, ist nicht nur die Anzahl der Datensätze, sondern insbesondere die Anzahl der Spalten ein ausschlaggebender Faktor für die Problemkomplexität.

Größe der Abhängigkeiten: Eine Möglichkeit, den exponentiellen Suchraum zu beschränken, ist die Größe der Abhängigkeiten zu limitieren, also die Anzahl der beteiligten Attribute. Zum Beispiel könnte man festlegen, dass Schlüsselkandidaten mit mehr als zehn Attributen irrelevant sind. Andererseits kann aber die vollständige Metadaten-Menge nützlich sein, um beispielsweise eine Relation auf der Basis ihrer funktionalen Anhängigkeiten zu normalisieren [18].

Anzahl der Abhängigkeiten: Während die meisten abhängigkeitsfokussierten Arbeiten, wie etwa Normalisierungstheorie oder das Schließen über Abhängigkeiten, von einer Handvoll an Abhängigkeiten als Input ausgeht, haben wir es üblicherweise mit Tausenden, Millionen oder sogar Milliarden an Abhängigkeiten in echten Datenbanken zu tun. So wird schon deren bloße Speicherung zu einem Problem, ganz zu schweigen von logischem Schlussfolgern oder einer Interpretation durch einen Experten.

Behandlung von Nullwerten: Die Semantik fehlender Werte ist ein spannendes Problem für jede Datenmanagement- und Analyseaufgabe, was auch auf Data Profiling zutrifft [11].

Komplexes Pruning: Huhtala et al. haben schon vor Längerem recht komplexe Regeln aufgezeigt, die den Suchraum für die Entdeckung funktionaler Abhängigkeiten beschneiden [10]. Wenn man Profiling für verschiedene Arten von Abhängigkeiten betreibt, wird es möglich, auch über Abhängigkeiten hinweg den Suchraum zu beschneiden.

1
2
3
4 **Relaxierte Abhängigkeiten:** Abgesehen von
5 strikten Abhängigkeiten ist es auch von
6 Interesse, partielle Abhängigkeiten zu
7 entdecken, also solche, die nur für einem
8 bestimmten Teil der Datenbank wahr sind, und
9 bedingte Abhängigkeiten (*Conditional*
10 *Dependencies*), die in einem wohl-definierten
11 Bereich wahr sind.

12 **Dynamische Daten:** Obwohl unser Fokus auf
13 Algorithmen für statische Datenmengen liegt,
14 sind wir auch daran interessiert, Metadaten für
15 sich verändernde Daten effizient aktuell zu
16 halten, ohne stets alles erneut zu ermitteln.

17 **Experimente:** Das Prüfen der Korrektheit von
18 Algorithmen für gegebene Echtweltdaten ist
19 recht unkompliziert. Hingegen gestaltet sich
20 das Generieren synthetischer Testdaten mit
21 speziellen Eigenschaften, wie zum Beispiel eine
22 bestimmte Anzahl und Verteilung an
23 funktionalen Abhängigkeiten, sehr schwierig.

24 **Ergebnisse interpretieren:** Entdeckte
25 Metadaten können nur für die aktuelle Instanz
26 validiert werden. Manche Abhängigkeiten
27 werden auch im Allgemeinen zutreffen, andere
28 wiederum nicht. Auf diese wohl wichtigste und
29 dabei schwierigste Herausforderung der
30 Interpretation und Nutzung von Data-Profiling-
31 Ergebnissen gehen wir in Abschnitt 2.4 ein.

32 Abschließend kann man feststellen, dass Data
33 Profiling weiter viele offene Forschungsfragen
34 anbietet!

35 2.2 Das Metanome-Projekt

36 Metanome ist unser Java-gestütztes
37 Framework und Werkzeug um relationale
38 Datensätze und Data-Profiling-Algorithmen zu
39 verwalten [16]. Unsere Motivation ist es, die
40 vielen in unserer Gruppe entwickelten und aus
41 anderen Arbeiten nachimplementierten
42 Algorithmen zu bündeln und eine einfache
43 Schnittstelle und Testumgebung für Entwickler
44 neuer Algorithmen zu schaffen. So ermöglichen
45 wir letztlich faire Vergleiche unter
46 konkurrierenden Algorithmen. Unser erstes
47 Augenmerk galt der Entdeckung funktionaler
48 Abhängigkeiten – in Metanome sind
49 mittlerweile acht veröffentlichte FD-
50 Entdeckungsalgorithmen implementiert,
51 inklusive derer, die in [17] ausgewertet
52 werden. Es kommen noch acht weitere
53 Algorithmen für andere Profiling-Aufgaben
54 hinzu (www.metanome.de).

55 2.3 RDF Profiling

Von besonderem Interesse, da reichhaltig,
vielfältig und oft allgemein verständlich, sind
Datensätze aus dem Bereich Linked (Open)
Data. Um diese zu untersuchen, wenden wir
sowohl herkömmliche, als auch neuartige
Data-Mining-Technologien für Linked Data in
deren RDF-Darstellung als Subjekt-Prädikat-
Objekt-Tripel an. Beispielsweise erlaubt uns die
Entdeckung häufiger Kombinationen (*Frequent*
Itemsets) von Eigenschaften (Prädikaten) oder
Objekten im Kontext von Subjekten, Daten mit
fehlenden Tripeln anzureichern. Eine weitere
Konfiguration – häufige Subjekte im Kontext
von Eigenschaften – erlaubt das Clustering von
Entitäten. Des Weiteren haben wir Data-
Mining-Technologien für die Entdeckung von
bedingten Inklusionsabhängigkeiten eingesetzt
[14]. Der typische Umfang solcher Daten (ein
populärer Datensatz stammt aus der Billion-
Tripels-Challenge und enthält gegenwärtig über
drei Milliarden Tripel) erfordert besonders
Speicher-effiziente Algorithmen.

Die meisten Ergebnisse dieser Arbeit fließen in
unser webbasiertes Werkzeug ProLOD++ ein,
welches Methoden anbietet, um in RDF-Daten
Schlüsselkandidaten zu entdecken, Klassen-
und Eigenschaftsverteilungen zu erforschen,
häufige grafische Muster zu entdecken und
vieles mehr (siehe [2] und Abbildung 1).

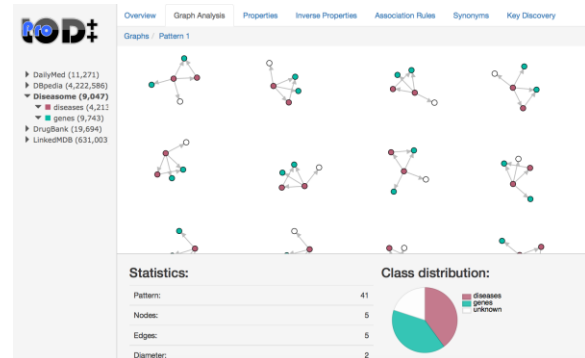


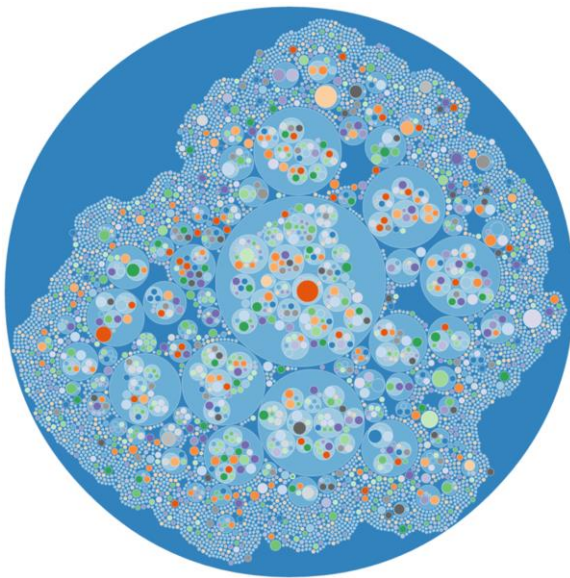
Abbildung 1: Erforschung häufiger Muster im
verknüpften Datensatz mit ProLOD++
(www.prolod.org)

56 2.4 Von Metadaten zu Semantik

Das Finden aller (und somit meist sehr vieler)
Abhängigkeiten in einer gegebenen
Datenmenge ist erst der Anfang einer
sinnvollen Analyse. Die überwiegende Mehrheit
an ermittelten Metadaten ist meist falsch: Sie
mögen nur in der aktuellen Instanz gültig sein,
oder aber durchaus in allen anderen möglichen
Instanzen gültig, aber dennoch bedeutungslos
sein. Die Spreu vom Weizen zu trennen ist sehr

1
2
3
4 schwierig, da es den Sprung von (Meta-)Daten
5 zu Semantik darstellt; nur ein Mensch kann
6 einen Schlüsselkandidaten zum Schlüssel
7 erheben, Inklusionsabhängigkeiten zu
8 Fremdschlüsseln machen oder funktionale
9 Abhängigkeiten zu einer Nebenbedingung
10 befördern.

11 Doch die Informatik kann helfen: Wir
12 verbringen derzeit einen Großteil unserer Zeit,
13 große Mengen an Metadaten in Schema-
14 Informationen umzuwandeln. Der erste Schritt
15 dabei ist die Entwicklung eines Metadaten-
16 Management-Systems, das Metadaten
17 speichert. Als nächstes entwickeln wir
18 Auswahl- und Rankingmethoden, um Nutzern
19 nur die vielversprechendsten Metadaten zu
20 präsentieren. Abschließend stellt die
21 Visualisierung von Metadaten ein wichtiges
22 Werkzeug dar, um Experten Hilfestellung zu
23 leisten und sie ihre Daten besser verstehen zu
24 lassen. Abbildung 2 zeigt beispielsweise
25 Zusammenhangskomponenten, die durch die
26 Entdeckung von Inklusionsabhängigkeiten in
27 Millionen von Web-Tabellen entstanden sind.



49 Abbildung 2: Cluster aus Web-Tabellen,
50 verbunden durch (sinnvolle)
51 Inklusionsabhängigkeiten

53 3. Datenbereinigung

54 Mit wachsenden Datenmengen wachsen auch
55 Probleme mit der Datenqualität. Eines der
56 interessantesten Probleme ist die mehrfache,
57 jedoch unterschiedliche Datenrepräsentation
58 des gleichen Objektes – so genannte
59 Duplikate. Diese Duplikate haben gleich

mehrere negative Auswirkungen: Zum Beispiel
können Bankkunden doppelte Identitäten
erhalten, Lagerbestände werden falsch
überwacht, Kataloge werden mehrfach an die
gleichen Haushalte geliefert, etc. Ein ähnliches
Problem ist das der Ähnlichkeitssuche in
strukturierten Daten: Zu einer Suchanfrage in
Form eines Datensatzes sollen die ähnlichsten
Datensätze in einer Datenbank gefunden
werden und es soll entschieden werden, ob
einer davon eine Übereinstimmung aufweist.

Beide Bereiche, Ähnlichkeitssuche und
Duplikaterkennung, erfahren zurzeit eine
Renaissance in Forschung und Industrie. Neben
der Erarbeitung wissenschaftlicher Beiträge
kooperieren wir mit Unternehmen zum
Technologietransfer. Sowohl unsere
Ähnlichkeitssuche als auch unsere
Duplikaterkennungstechniken werden aktiv von
Industriepartnern genutzt.

3.1 Duplikaterkennung

Das Aufspüren von Duplikaten in einer
gegebenen Tabelle ist schwierig: Erstens sind
die Darstellungen der Duplikate normalerweise
eben nicht identisch, sondern deren Werte
weichen etwas voneinander ab. Zweitens
müssten im Allgemeinen alle Datensatzpaare
verglichen werden, was für große
Datenmengen undurchführbar ist. Unsere
Forschung untersucht beide Aspekte, erstens
durch den Entwurf effektiver Ähnlichkeitsmaße
und zweitens durch das Entwickeln effizienter
Algorithmen zur Reduzierung des Suchraumes.

Ein Fokus unserer Arbeit ist es, verbesserte
Variationen der eleganten und einfachen
Sorted Neighborhood-Methode zu entwickeln
[9], zum Beispiel durch Anpassung an
verschachtelte XML-Daten, durch progressive
Verarbeitung, die möglichst früh möglichst
viele Ergebnisse liefert, durch Parallelisierung
für GPU-Verarbeitung oder durch eine adaptive
Version, die beweisbar effizienter ist als das
Original [5].

Aus unserer Erfahrung heraus leiden die
meisten Projekte zum Thema
Duplikaterkennung bei dem Versuch,
Technologie und Methodik in die industrielle
Wirklichkeit zu transferieren. Die Verfügbarkeit
von Daten ist ein erstes Problem, welches
aufkommt, selbst wenn eine Kooperation fest
etabliert ist und alle teilnehmenden Parteien
grundsätzlich dem Projekt zustimmen. Als
nächstes werden Domänen- und Partner-
spezifische Ähnlichkeitsmaße benötigt, die zum

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 jeweiligen Anwendungsfall passen.
5 Unternehmen haben oftmals sehr
6 unterschiedliche Sichtweisen darauf, was ein
7 Duplikat überhaupt ist: Das Messen von *Recall*,
8 also der Vollständigkeit der entdeckten
9 Duplikate, ist unmöglich, in Anbetracht der
10 großen Datenmengen und *Precision*, also der
11 Korrektheit der entdeckten Duplikate, ist
12 überraschend dehnbar, je nachdem, wen man
13 zur Validierung befragt. Abschließend ist zu
14 sagen, dass die reale Welt viele praktische
15 Details vorhält, die in der Forschung bequem
16 ignoriert werden können. Mit [20] gelang es
17 uns, solche Schwierigkeiten zu überwinden und
18 die Datenqualität unserer Partner zu
19 verbessern.

20 **3.2 Ähnlichkeitssuche**

21 Ein Problem, das mit der Duplikatsuche zwar
22 verwandt ist, jedoch abweichende
23 Anforderungen hat, ist die effiziente
24 Ähnlichkeitssuche. Anstatt alle oder zumindest
25 sehr viele Datensatzpaare offline zu
26 vergleichen ($n \times n$), fragt die Ähnlichkeitssuche
27 online nach allen Datensätzen, die zu einem
28 gegebenen Anfragedatensatz passen ($1 \times n$).
29 Ein typisches Fallbeispiel ist ein Callcenter-
30 Mitarbeiter, der Kundendaten nur auf Basis des
31 Kundennamens und Wohn-orts aus dem
32 System auslesen möchte. Die größte
33 Herausforderung ist dabei die Entwicklung
34 eines passenden Ähnlichkeits-Index, eine
35 weitaus komplexere Datenstruktur als die
36 üblichen Indizes, die auf Gleichheit von Werten
37 beruhen.

38
39 Eine unserer Lösungen basiert auf klassischen
40 Ansätzen der Anfrageoptimierung: Wir wählen
41 Ähnlichkeitsindexzugriffe aus, basierend auf
42 ihrer Selektivität und ihrer Kosten, die jeweils
43 durch einen dynamisch gewählten Schwellwert
44 modifiziert werden: Ein niedriger Schwellwert
45 ergibt mehr Kandidaten, bedeutet aber auch
46 mehr Zugriffe und somit höhere Kosten, um
47 die Kandidaten abzurufen und abschließend zu
48 vergleichen [15]. Eine weitere Erkenntnis ist
49 die Bedeutung von Häufigkeiten für
50 Ähnlichkeitsmaße: Je nach Häufigkeit der
51 Anfragewerte sollten unterschiedliche
52 Gewichtungen vorgenommen werden
53 (Leutheusser-Schnarrenberger vs. Müller).

54
55 Derzeit erweitern wir diese Ideen, um das
56 Problem eines ständig wachsenden
57 Datensatzes zu lösen, der Duplikat-frei
58 gehalten werden soll: Jede Anfrage ist zugleich
59 ein Einfügen in den Datenbestand.

4. TEXT MINING

Unstrukturierte Daten in Form von
natürlichsprachlichen Textdokumenten findet
man überall, wo Kommunikation und
Informationsaustausch zwischen Menschen
stattfindet. Entsprechend vielfältig sind die
Dokumentarten. Insbesondere der Erfolg des
Internets als Kommunikationsmedium hat zu
einer Explosion öffentlich zugänglicher
Textsammlungen geführt. Neben
benutzergeneriertem Inhalt, wie beispielsweise
Wikipedia, Blogs oder Tweets, gibt es auch
eine Menge professionell erstellter Inhalte, wie
zum Beispiel Patente, Zeitungsartikel,
politische Reden, Romane oder
wissenschaftliche Publikationen. Weniger frei
zugänglich, aber auch in großen Mengen
vorhanden, sind natürlichsprachliche
Dokumente wie Patientenakten, Geschäfts-
Emails oder Instant-Messaging-Nachrichten.
Für das Text Mining stellen diese sehr
heterogenen Daten eine besondere
Herausforderung dar, da jedes Genre seine
eigenen Charakteristiken besitzt, und
Dokumente von wenigen Worten bis tausende
von Seiten lang sein können. Um diese
vielfältigen Daten zu analysieren, Muster zu
erkennen und Wissen zu generieren, bedient
sich das Text Mining klassischer Methoden des
Information Retrievals, der automatische
Sprachverarbeitung, der
Informationsextraktion und des maschinellen
Lernens. Aktuelle Projekte am Fachgebiet sind
die Analyse von Nachrichten und von
Geschäftskommunikation, die Erforschung und
Weiterentwicklung von Topic-Modellen, die
gemeinsame Analyse mehrerer
unterschiedlicher Textsammlungen, sowie die
Erforschung des zeitlichen Aspekts bei
dynamischen, sich schnell ändernden
Textsammlungen.

4.1 Analyse von Nachrichten

Nachrichten dienen als Grundlage mehrerer
Text-Mining-Aufgaben, wie beispielsweise der
automatischen Erzeugung von
Zusammenfassungen oder dem Erkennen von
Ereignissen. Uns interessieren Zeitungsartikel
im Zusammenhang mit Medien-Bias. Um
diesen zu erkennen, genügt die Analyse eines
einzelnen Zeitungsartikels nicht, und meist
auch nicht ein isolierter Blick auf alle Artikel
einer Zeitung. Wir bedienen uns daher
zusätzlicher Informationsquellen, und zwar
parlamentarischer Reden sowie
Leserkomentaren.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Vergleiche von parlamentarischen Reden mit Nachrichtenartikeln verschiedener deutscher Nachrichtenagenturen haben in ersten Experimenten gezeigt, dass der wahrgenommene Bias automatisch quantifiziert werden kann [12]. Das Benutzen eines bestimmten Vokabulars ist ein erster Indikator für Bias (z.B. „Kernenergie“ vs. „Atomenergie“ in Deutschland). Neben bestimmten Standpunkten (Statement Bias) können Zeitungen ihre Leser beeinflussen, indem sie nur über bestimmte Themen berichten (Gate-Keeping Bias) oder indem sie bestimmte Positionen stärker als andere abdecken (Coverage Bias). Die automatische Erkennung aller drei Bias-Arten ist besonders schwierig, da in der deutschen Medienlandschaft nur subtile Unterschiede zwischen den Leitmedien bestehen. Das macht es notwendig, nicht nur Entitäten (Politiker, Parteien, Experten) und ihre Beziehungen zu extrahieren und zu analysieren, sondern ebenso detaillierte Opinion-Mining und Sentiment-Analyse zu betreiben und alles ins Verhältnis zu einer virtuellen, objektiven Berichterstattung zu setzen.

Erste Ergebnisse konnten wir auch durch die Analyse von Zeitungsartikelkommentaren erzielen. Unsere Hoffnung war, dass Leser in ihren Kommentaren den Bias der Zeitung unverblümter offenbaren [6], was wir auch anhand einer Klassifizierungsaufgabe bestätigen konnten. Ein weiterer Schwerpunkt liegt auf dem Erkennen von Hasskommentaren, um im Idealfall die Veröffentlichung ebendieser automatisch zu unterbinden. Dafür arbeiten wir mit einer großen deutschen Tageszeitung zusammen.

4.2 Analyse von Geschäftskommunikation

Das Erkennen von Entitäten, insbesondere Firmennamen, und Themen in E-Mails ist der erste Schritt, um komplexe Zusammenhänge innerhalb eines Unternehmens und dessen Kunden zu erkennen. Im Rahmen eines Industrieprojektes mit einer großen deutschen Bank streben wir den Aufbau von Unternehmensnetzwerken zur Unterstützung ihrer Risikomanagementabteilung an. Diese Firmennetzwerke werden automatisch aus Zeitungsartikeln extrahiert und stellen eine neue Herausforderung für die benannte Entitätserkennungsaufgabe (Named Entity Recognition) dar, die aufgrund komplexer, oft zweideutiger Benennung besonders schwierig für deutsche Firmennamen ist. Darüber hinaus unterscheiden sich die Beziehungstypen, die

für uns interessant sind, von üblichen, binären Beziehungen, wie „verheiratet mit“ oder „wohnhaft in“. Unsere Unternehmensnetzwerke erfordern die Erfassung von Beziehungen, die nicht notwendigerweise binär sein müssen, z.B. „Wettbewerber mit“ oder „Lieferant für“. Zu diesem Zweck haben wir einen Algorithmus entwickelt, welcher anhand weniger vorgegebener Instanziierungen einer Relation andere, gleichartige Beziehungen findet. Der Algorithmus basiert auf Snowball [3] und kann mit jeder Art von Beziehung, die vom Benutzer bereitgestellt wird, umgehen.

4.3 Topic-Modelle

Topic-Modelle sind statistische Modelle, welche in großen Dokumentsammlungen Wörter thematisch in sogenannten Topics gruppieren und damit einen Überblick über die vorhandenen Themen in einer Textsammlung liefern. Neben dem Einsatz dieser Modelle für unsere Analysen arbeiten wir auch an der Weiterentwicklung und Anpassung dieser Modelle für diverse Anwendungsszenarien. So haben wir beispielsweise die teilweise vorhandenen Schlagwörter bei Projektanträgen benutzt, um Projekte in Themenbereiche einzuordnen [19]. Dies hat auch die Möglichkeit eröffnet, die Summen der geförderten Projekte auf Themen umzulegen und Trends in der Gesundheitsforschungsförderung offenzulegen.

Ein weiterer Schwerpunkt liegt auf dem Einbinden sogenannter Word Embeddings zur Verbesserung von Topic-Modellen. Word Embeddings erlauben die Repräsentation von Wörtern in n-dimensionalen, reellen Vektorräumen und somit das Finden von semantisch ähnlichen Wörtern. Diese Technologie, welche auf Deep Learning beruht, machen wir uns zunutze, um allgemeine Wörter durch thematische relevantere Wörter in den Topics zu ersetzen.

4.4 Analyse mehrerer Textsammlungen

Ein besonderes theoretisches Problem, mit dem wir uns beschäftigen, ist die Analyse von Textdaten über Korpusgrenzen hinweg, was schon bei Zeitungsartikeln und Parlamentsreden nötig war. Hier haben wir eine Reihe praktischer Anwendungen untersucht, um systematisch Schwierigkeiten zu erkunden.

In einem Versuch, die Kluft zwischen traditionellen Nachrichten und Social Media zu

1
2
3
4 überbrücken, haben wir ein Tweet-
5 Empfehlungs-System entwickelt [13]. Das Ziel
6 war es, dem Leser eines News-Artikels über ein
7 bestimmtes Ereignis einen Überblick über
8 Reaktionen auf Twitter zu geben. Während
9 Twitter häufig genutzt wird, um Informationen
10 zu teilen und zu verbreiten, wird es ebenso
11 häufig genutzt, um Meinungen auszudrücken,
12 Ideen abzulehnen oder bestimmte Standpunkte
13 zu unterstützen. Um diese Meinungen zu
14 erkennen, müssen traditionelle
15 Sentimentanalysetechniken angepasst werden,
16 um Emojis, Abkürzungen, Slang usw. zu
17 erkennen. Das Missverhältnis zwischen der
18 Sprache, die in Nachrichtenartikeln und Tweets
19 verwendet wird, macht die Empfehlung des
20 einen aufgrund des anderen sehr schwierig.

21 **4.5 Analyse dynamischer Korpora**

22
23 Viele der analysierten Textsammlungen sind
24 nicht statisch, sondern ändern sich mehr oder
25 weniger schnell. Die Einbeziehung der
26 zeitlichen Komponente bei der Analyse, aber
27 auch bei konkreten Aufgaben, wie dem
28 Empfehlen von Objekten, muss deshalb
29 besonders Rechnung getragen werden. Twitter
30 stellt durch die kurzen Textbeiträge und der
31 hohen Frequenz an neuen Tweets eine
32 besondere Herausforderung dar. Durch
33 verschiedene zeitliche Modelle haben wir
34 versucht, die Dynamik von Hashtags zu
35 analysieren und die Ergebnisse zu verwenden,
36 um das Empfehlen von Hashtags zu verbessern
37 [8].

38
39 Weiter haben wir die längerfristigen,
40 thematischen Veränderungen in einem Forum
41 untersucht [7]. Hierbei konnten nicht nur
42 unterschiedliche Softwareversionen anhand der
43 zeitlichen Diskussion identifiziert werden,
44 sondern auch der Wandel bei der Benutzung
45 von Begriffen und der sich ändernde Kontext
46 für bestimmte Wörter.

47 **5. Lehre am Fachgebiet**

48
49 Im Bachelorstudium bieten wir, neben den
50 typischen Vorlesungen zu den Grundlagen und
51 der Implementierung von
52 Datenbankmanagementsystemen und den
53 sporadischen Proseminaren, jährlich ein bis
54 zwei Bachelorprojekte an – eine besondere
55 Veranstaltungsform am HPI. Zwischen vier und
56 acht Bachelorstudenten bearbeiten zum
57 Abschluss ihres Studiums ein Softwareprojekt
58 über einen Zeitraum von zwei Semestern und
59 stets in Kooperation mit einem externen
60 Partner, der die Aufgabenstellung vorgibt und

gleichsam als „Kunde“ das Projekt aktiv
begleitet. Die individuellen Bachelorarbeiten
greifen Themen des Bachelorprojekts auf. Die
Studenten sind im Wintersemester circa zur
Hälfte ihrer Zeit im Rahmen des Projekts tätig,
im darauf folgenden Sommersemester
bearbeiteten sie in der Regel ausschließlich die
Aufgaben im Projekt. Aus der Studienordnung:
„Es handelt sich um praxisnahe Projekte, bei
denen die Studierenden nicht nur als
Entwickler kreativ werden, sondern in denen
sie auch die besonderen Merkmale der
Koordination von vielen Projektbeteiligten
erleben.“ Die entstandenen Softwarelösungen
werden zum Semesterabschluss der
Öffentlichkeit vorgestellt und nicht selten
anschließend durch die Partnerorganisationen
produktiv eingesetzt. Unsere bisherigen
Projektpartner waren u.a. Wikimedia
Deutschland, Commerzbank, IBM, Capgemini
sowie diverse kleine und mittelständige
Unternehmen
(<https://hpi.de/naumann/teaching/bachelorprojekte.html>).

Im Masterstudium bieten wir Vorlesungen und
Seminare an, deren Themen durch unsere
Forschungsrichtungen bestimmt werden. Zu
den wichtigsten und wiederkehrenden
Vorlesungen gehören „Information
Integration“, „Data Profiling“, „Data Mining and
Probabilistic Reasoning“ sowie „Information
Retrieval and Web Search“. Seminare ergänzen
den Stoff und werden des Öfteren als
Projektseminare durchgeführt, bei denen
neben den üblichen Vorträgen und
Ausarbeitungen auch vergleichende
Implementierungen erstellt werden. Auch im
Masterstudium ist eine Teamarbeit vorgesehen,
nämlich ein Masterprojekt mit drei bis sechs
Studenten, die im Gegensatz zu den
Bachelorprojekten eine konkrete
Forschungsfrage des Fachgebiets untersuchen
sollen. Ein typisches Ergebnis eines solchen
Projekts ist die (oft erfolgreiche) Einreichung
eines Manuskripts auf einer wissenschaftlichen
Konferenz
(<https://hpi.de/naumann/teaching/masterprojects.html>).

Danksagung Unserer Forschung genoss die
Unterstützung verschiedener Partner wie der DFG
und Unternehmen, die sich für das Verständnis und
die Verbesserung ihrer Daten interessieren. Die hier
vorgestellten Arbeiten beruhen – natürlich – auf der
Forschung unserer hervorragenden Doktoranden:
Tobias Bleifuß, Toni Grütze, Hazar Harmouch,
Maximilian Jenders, Anja Jentsch, John Koumarelas,
Sebastian Kruse, Konstantina Lazaridou, Michael

1
2
3
4 Loster, Thorsten Papenbrock, Julian Risch, Ahmad
5 Samiei und Zhe Zuo.

6
7 Zwei weitere HPI-Fachgebiete, mit denen wir
8 kooperieren, arbeiten ebenfalls in der Datenbank-
9 Community: Das Fachgebiet „Enterprise Platforms
10 and Integration Concepts“ (EPIC) unter der Leitung
11 von Hasso Plattner und Matthias Uflacker sowie das
12 Fachgebiet „Knowledge Discovery und Data Mining“
13 (KDD) unter der Leitung von Emmanuel Müller.

14 **Literatur**

- 15 1. Abedjan Z, Golab L Naumann F (2015)
16 Profiling relational data: a survey. VLDB
17 Journal, 24(4):557-581
- 18 2. Abedjan Z, Gruetze T, Jentzsch A, Nau-
19 mann F (2014) Profiling and mining RDF
20 data with ProLOD++. In Proc. of the Inter-
21 national Conference on Data Engineering
22 (ICDE), Demo. S 1198-1201
- 23 3. Agichtein E, Gravano L (2000) Snowball:
24 Extracting relations from large plain-text
25 collections. In Proc. of the ACM Conference
26 on Digital Libraries. S 85-94
- 27 4. Agrawal D, Bernstein P, Bertino E,
28 Davidson S, Dayal U, Franklin M, . . .
29 Widom J (2012) Challenges and opportuni-
30 ties with Big Data. Technical report, Com-
31 puting Community Consortium,
32 [http://cra.org/ccc/docs/init/bigdatawhitepa-
33 per.pdf](http://cra.org/ccc/docs/init/bigdatawhitepaper.pdf)
- 34 5. Draisbach U, Naumann F, Szott S, Wonne-
35 berg O (2012) Adaptive windows for dupli-
36 cate detection. In Proc. of the International
37 Conference on Data Engineering (ICDE). S
38 1073-1083
- 39 6. Godde C, Lazaridou K, Krestel R (2016)
40 Classification of German Newspaper Com-
41 ments. In Proc. of the Conference Lernen,
42 Wissen, Daten, Analysen (LWDA). S 299-
43 310
- 44 7. Gruetze T, Krestel R, Naumann F (2016)
45 Topic Shifts in StackOverflow: Ask it like
46 Socrates. In Proc. of the 21st International
47 Conference on Applications of Natural Lan-
48 guage to Information Systems (NLDB),
49 volume 9612. S 213-221
- 50 8. Gruetze T, Yao G, Krestel R (2015) Learn-
51 ing temporal tagging behaviour. In Proc. of
52 the Temporal Web Analytics Workshop
53 (TempWeb) at the International World
54 Wide Web Conference (WWW). S 1333-
55 1338
- 56 9. Hernández MA, Stolfo SJ (1998) Real-world
57 data is dirty: Data cleansing and the
58 merge/purge problem. Data Mining and
59 Knowledge Discovery, 2(1):9-37
- 60 10. Huhtala Y, Kärkkäinen J, Porkka P, Toivo-
61 nen H (1999) TANE: An efficient algorithm

for discovering functional and approximate
dependencies. Computer Journal,
42(2):100-111

11. Köhler H, Link S, Zhou X (2015) Possible
and certain SQL keys. In Proc. of the VLDB
Endowment, 8(11):1118-1129
12. Krestel R, Wall A, Nejd W (2012) Treehug-
ger or Petrolhead? Identifying Bias by
Comparing Online News Articles with Politi-
cal Speeches. In Proc. of the International
World Wide Web Conference (WWW). S
547-548
13. Krestel R, Werkmeister T, Wiradarma TP,
Kasneji G (2015) Tweet-recommender:
Finding relevant tweets for news articles.
In Proc. of the International World Wide
Web Conference (WWW). S 53-54
14. Kruse S, Jentzsch A, Papenbrock T, Kaoudi
Z, Quiane-Ruiz JA, Naumann F (2016)
RDFind: Scalable conditional inclusion de-
pendency discovery in RDF datasets. In
Proc. of the International Conference on
Management of Data (SIGMOD). S 953-
967
15. Lange D, Naumann F (2011) Efficient simi-
larity search: Arbitrary similarity measures,
arbitrary composition. In Proc of the Inter-
national Conference on Information and
Knowledge Management (CIKM). S 1679-
1688
16. Papenbrock T, Bergmann T, Finke M,
Zwiener J, Naumann F (2015) Data
profiling with Metanome (demo). In Proc.
of the VLDB Endowment, 8(12):1860-1871
17. Papenbrock T, Ehrlich J, Marten J, Neubert
T, Rudolph JP, Schönberg M, Zwiener J,
Naumann F (2015) Functional dependency
discovery: An experimental evaluation of
seven algorithms. Proceedings of the VLDB
Endowment, 8(10):1082-1093
18. Papenbrock T, Naumann F (2017) A Hybrid
Approach for Efficient Unique Column
Combination Discovery. In Proc. der
Fachtagung Business, Technologie und
Web (BTW), accepted
19. Park J, Blume-Kohout M, Krestel R,
Nalisnick E, Smyth P (2016) Analyzing NIH
Funding Patterns over Time with Statistical
Text Analysis. In Scholarly Big Data: AI
Perspectives, Challenges, and Ideas, Work-
shop at AAAI. S 698-704
20. Weis M, Naumann F, Jehle U, Lufter J,
Schuster H (2008) Industry-scale duplicate
detection. In Proc. of the VLDB Endow-
ment, 1(2):1253-1264
21. Zuo Z, Kasneji G, Gruetze T, Naumann F
(2014) BEL: Bagging for entity linking. In
Proc. of the International Conference on

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

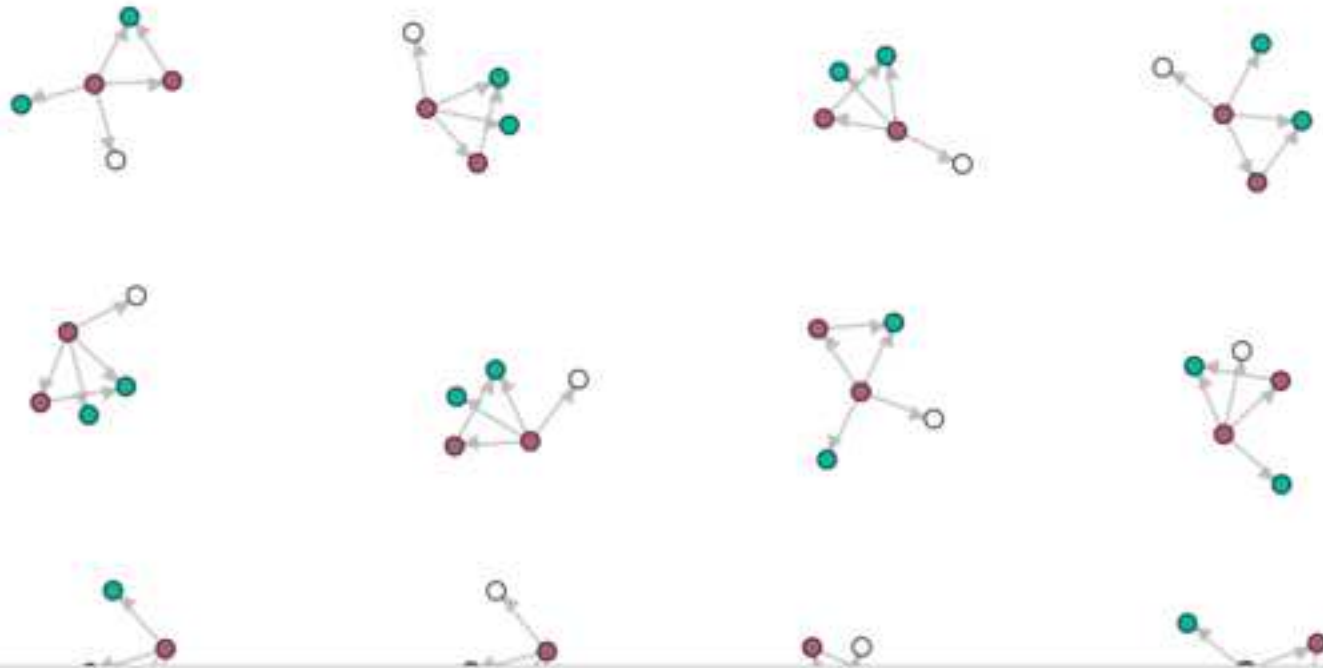
Computational Linguistics (COLING). S
2075-2086



Overview Graph Analysis Properties Inverse Properties Association Rules Synonyms Key Discovery

Graphs / Pattern 1

- ▶ DailyMed (11,271)
- ▶ DBpedia (4,222,586)
- ▼ **Diseasome (9,047)**
 - ▾ diseases (4,213)
 - ▾ genes (9,743)
- ▶ DrugBank (19,694)
- ▶ LinkedMDB (631,003)



Statistics:

Pattern:	41
Nodes:	5
Edges:	5
Diameter:	2

Class distribution:

